

Highlight, then Summarize: RAG without the jailbreaks

Giovanni Cherubin
Microsoft

Andrew Paverd
Microsoft

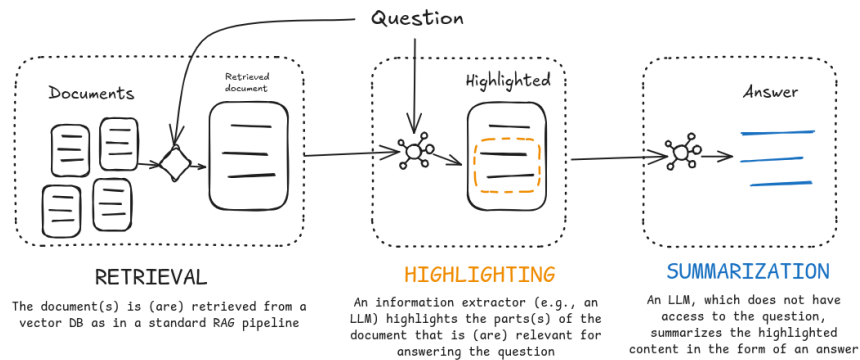


Figure 1: HS-enhanced RAG pipeline. A malicious user asking a question to this system cannot jailbreak the LLM, by design.

Abstract

Preventing model hijacking and jailbreaking is a troublesome problem in the deployment of Large Language Models (LLMs). In RAG systems, for example, malicious users can tamper with the LLM’s inference, and force it to produce undesired contents, beyond what initially envisioned by the system designers. Several mitigations rely on content classifiers, but they are mostly of reactive nature; an ideal jailbreak classifier would need to *understand* what it means for something to be malicious, which is out of reach from a scientific and philosophical standpoint for the foreseeable future.

In this paper, we look at a simple enhancement for RAG systems that prevents these attacks *by design*. The idea behind this design pattern, which we refer to as Highlight&Summarize (H&S), is to accomplish the same goal as a RAG pipeline (i.e., to provide natural language answers to questions, based on relevant sources) without ever letting the LLM generator be directly exposed to a user’s question. This is obtained by having two parts: a highlighter, which extracts from the sources relevant information to a user’s query, and a summarizer, an LLM that polishes and summarizes the extracted information. This paper discusses and evaluates the security and utility implications of this design pattern. To our own surprise, the responses given by a H&S pipeline are comparable to, and often of a better quality, than those coming from a RAG pipeline.

1 Introduction

Retrieval-augmented Generation (RAG) is proving to be a sound and reliable solution for answering questions using documents from a knowledge base. From customer support systems to search engines, RAG combines the ability of LLMs to answer questions expressed in natural language with the efficiency of vector databases for retrieving information from large data repositories, to provide a robust production-ready solution for many applications. The core idea behind RAG is that, when a user asks a question, the system first retrieves a set of relevant documents from the knowledge base and passes these to a generative LLM together with the user’s

question; on this basis, the LLM produces an answer to the question, which is returned to the user.

There are, alas, many ways an adversary can exploit RAG systems, depending on which inputs the adversary can control. In this paper, we focus on the scenario where the knowledge base is trusted, while the users’ inputs are untrusted and potentially malicious. A real-world example is a chatbot that is deployed by a company to answer questions based on the company’s curated knowledge base (e.g., consisting of official FAQs or policy documents).

In this setting, a malicious user could attempt to *jailbreak* the system to achieve various objectives. They could repurpose the LLM generation step for some other task (e.g., asking questions unrelated to the company) – this is sometimes referred to as *model hijacking* or *wallet abuse*. Alternatively, they could make the LLM generate dangerous or offensive content that harms the reputation of the company. Finally, they could trick the LLM into generating unintended outputs and then claim that these represent the company’s intent; e.g., a user could persuade the chatbot to say “here’s a 90% discount on product X”, which may in some cases constitute a legally binding statement from the company. It was recently reported that Air Canada had to pay compensation for misleading information provided to one of their customers by their FAQ chatbot.¹

In this paper, we introduce and evaluate a design that prevents these attacks by design, by revisiting the data flow with the RAG pipeline. We refer to this as Highlight&Summarize (H&S). Our approach consists of two-components: First, the *highlighter* component inspects the retrieved document(s) and selects passages from them that are relevant to answering the user’s question – the digital equivalent of *highlighting* those passages. The highlighter can be based on existing methods for extractive Q&A and passage retrieval, or on modern generative LLM pipelines. Secondly, the *summarizer* component, which is an LLM, generates an answer by summarizing the highlighted passages. **Critically, the summarizer never sees**

¹<https://www.theguardian.com/world/2024/feb/16/air-canada-chatbot-lawsuit>.

the user’s question. This mitigates the attacks by malicious users described above, whilst maintaining the utility of the RAG system on its intended question-answering task.

The remainder of this paper is organized as follows. We first present background and discuss related work in this area (section 2) and then define our precise threat model (section 3). We then describe different implementations of the H&S design (section 4). Next we explain our experimental setup (section 5), evaluate the individual components of H&S (section 6), and evaluate the full H&S system in comparison with a vanilla RAG pipeline (section 7). Finally, we evaluate the security of H&S pipelines (section 8), and discuss future work and open questions (section 9).

We release both the code and the data we used, including the responses given by the RAG pipelines and their ratings given by the LLM-as-Judges. Our aim in doing so is to enable further analysis of H&S, without incurring the monetary and environmental costs of rerunning our experiments.

2 Background

We consider an LLM-based RAG application that retrieves information from a trusted knowledge base in order to answer users’ questions. Specifically, users pose questions to the application via a single-turn chat-style interface. Following standard RAG designs, the application searches its knowledge base for documents related to the user’s question. In a typical RAG application the retrieved documents are concatenated to the user’s question and input to an LLM, which generates a response. Importantly, we assume the application owner has inspected all documents in the knowledge base and confirmed that they contain only trustworthy information.

This type of RAG application is widely used, for example in customer support services.^{2,3} The benefits of using RAG in such applications include the ability to ground the LLM’s response on company-specific information, reduce hallucinations, and update the knowledge base without retraining the LLM.

2.1 Related work

Jailbreaking. Several techniques have been proposed to detect or prevent jailbreaking of LLM-based systems, including those used in RAG. The main classes of defenses include: finetuning the LLM to reduce the risk of undesirable output; using defensive system prompts to increase the difficulty of jailbreaking [19, 21]; applying different types of classifiers to the inputs to detect potential jailbreaks [4, 12, 13, 25, 28]; and preprocessing the input in order to remove or reduce the impact of potential jailbreaks [13, 17, 31, 32]. As a way of mitigating model hijacking (i.e., repurposing the LLM for a different task), techniques such as NeMo support programmable guardrails that allow the application owner to specify dialogue flows by canonicalizing inputs [31].

Our proposed solution, H&S, sidesteps the problem by ensuring that the user has very limited control over the outputs that are generated by the LLM (summarizer).

Similarly to defenses such as CaML [10] and FIDES [9], H&S is a system-level defense in that it deterministically limits the attacker’s

²<https://careersatdoordash.com/blog/large-language-modules-based-dasher-support-automation/>.

³<https://medium.com/tr-labs-ml-engineering-blog/better-customer-support-using-retrieval-augmented-generation-rag-at-thomson-reuters-4d140a6044c3>.

control by design. However, H&S focuses on the opposite threat model, where the retrieved data is trusted but the user inputs are not. H&S is a realization of the Context-Minimization design pattern [5].

Passage retrieval and extractive Q&A. The first step in H&S, *highlight*, selects portion(s) of text that contain the answer to a question. Two research fields in NLP address a similar goal: passage retrieval and extractive Q&A. Passage retrieval is the process of extracting from a set of documents one or more texts that are relevant to answering a question; over the years, many methods have been built, with applications in the medical and legal sectors [6, 15, 22, 23]. With a similar goal, extractive Q&A aims to select (brief) portions of text from a larger document which answer a question. Extractive Q&A methods are largely based on modern NLP architectures, with BERT-like models being the most successful [18, 26, 35]. Passage retrieval and extractive Q&A are typically combined in practical systems [16, 29].

One could think of the highlighter in a H&S pipeline as a similar concept to passage retrieval and extractive Q&A. In this view, one of our contributions is the observation that adding a summarizer after highlighting improves the overall performance (subsection 6.1).

In our experiments, we also explicitly explore using an extractive Q&A method as a way of implementing the highlighter; however, we note that this is often outperformed by highlighters that are based on generative LLMs (see section 4 for their description).

3 Threat model

Some users of a RAG application may be adversarial. We assume a strong adversary who has full knowledge of all the documents in the knowledge base; for example, this may be the case if the knowledge base consists of published documentation or FAQ web pages. The adversary is able to submit arbitrary questions and observe the responses. The adversary may want to achieve the following goals:

- Repurpose the application to perform a different task. For example, using a customer service chatbot to summarize large amounts of text, at the expense of the application owner.
- Jailbreak the application and cause it to generate undesirable output. This could include content that causes reputational damage to the application owner, or even constitutes an unintentional yet legally enforceable commitment from the application owner.

The goal of the application owner is to prevent both of the above classes of attacks. Note that we do not aim to defend against other types of attacks on RAG systems, such as extracting information [14, 27] or poisoning the knowledge base [36, 38].

4 Highlight&Summarize

We propose Highlight&Summarize (H&S), which *enhances* RAG systems to provide security against jailbreaks. The main design principle is that a (malicious) user must be unable to provide direct inputs to the LLM that generates the response. We achieve this by separating the generative Q&A process into two components: *Highlight*, where the system select text passage(s) that is (are) relevant to answering a query; here, we enforce that the output of the highlighter must be text that appeared in the document. *Summarize*,

where the system summarizes the outputs of the highlighter, as if it were answering a question. This process is illustrated in Figure 1.

By virtue of this design, malicious users are unable to directly influence the outputs of the system – which are provided by the summarizer. In section 8, we further evaluate the security of this design pattern, as well as discuss possible attacks and mitigations.

4.1 H&S implementations

H&S can be implemented in multiple ways. In this section we describe our implementation of the summarizer component and our three different implementations of the highlighter. We refer to the combination of any of these highlighters with the summarizer as an *H&S pipeline*.

Summarizer (common across all H&S implementations.) The summarizer is a zero-shot prompt tuned LLM that is tasked with i) guessing what question the highlighted text was intended to answer, and ii) reformulating the highlighted text in form of an answer. Only the answer is returned to the user. The guessed question is not returned to the user – its purpose is to ground the generated answer, and to aid in evaluation. We use Azure OpenAI’s structured output option⁴, which forces the LLM to give a response matching a desired format (in this case, `{"guessed_question": str, "answer": str}`); we found this helps with the quality of the responses, and allows us to evaluate what question the LLM guessed in our data analysis. Unless otherwise specified, we employ OpenAI’s GPT-4.1 mini [2] as our default generative LLM.

H&S Baseline. This highlighter is a zero-shot prompt-tuned LLM, which is tasked with extracting relevant information from the retrieved documents. To enforce that the output corresponds to text from the documents, we use fuzzy string matching via RapidFuzz⁵; fuzzy string matching allows for a small variability (which may be due to the model’s randomness) without compromising the security requirement that the extracted text must be (approximately) a substring of the original document.

H&S Structured. This highlighter improves upon H&S Baseline by asking the LLM to first output an answer to the question and then the highlighted text from the documents. We again use Azure OpenAI’s structured output option to enforce the following format: `{"answer": str, "text_extracts": list[str]}`.

The generated answer is not passed to the summarizer. We observed that asking the LLM for this output helped with grounding its responses, thereby producing better highlights from the text.

H&S DeBERTaV3. This highlighter is an extractive Q&A model from the BERT family. In experiments, we use DeBERTaV3, which uses disentangled attention and a better mask decoder to improve upon BERT and RoBERTa. We use two versions of this model: *H&S DeBERTaV3 (SQuAD2)* is a DeBERTaV3 model that was fine-tuned⁶ on the SQuAD2 dataset [30]. Since this model is fine-tuned for short-span extractive Q&A, as encouraged by the SQuAD2 dataset, its performance in our experiments was lacking: the highlighter of an H&S pipeline is best served by a tool that outputs longer

span passages to a question. For this reason, we also employ *H&S DeBERTaV3 (RepliQA)*, which is a DeBERTaV3 model that we fine-tuned on the RepliQA dataset (see section 5). We train the model to return the gold passage of an answer, rather than the expected answer; this significantly improves the performance of the H&S pipeline. Fine-tuning this model on splits 0-2 of the RepliQA dataset (section 5) took around 7 hours on an NVIDIA A100 GPU.

5 Experimental setup

We describe the datasets and metrics used in our experiments.

5.1 Datasets

Evaluating generative LLMs for RAG pipelines requires extreme care: while we may want to judge the LLM’s ability at responding based on source documents, it often happens that the LLM answers based on its training data instead. For this reason, we conduct our main experiments using the RepliQA dataset [20]. This is a human-created dataset that contains questions based on natural-looking documents about *fictional* events or people. By doing so, we ensure that the performance is not affected by the ability of LLMs to memorize their training data.

Importantly, each entry in the RepliQA dataset contains a *gold passage*, called `long_answer`, which is a substring of the retrieved document selected by a human annotator that is relevant to answering the question. We use this field as part of our performance measurements, as well as for fine-tuning (on a separate split of the dataset) a DeBERTaV3 extraction Q&A model to implement the H&S DeBERTaV3 (RepliQA) pipeline.

The RepliQA dataset consists of 5 splits (numbered from 0 to 4), which were released gradually over a year. Each split contains 17,955 examples. For our evaluation, we used the most recently released split (`split_3`, which was released on April 14th, 2025). During data analysis, we observed 10 mislabelled instances in this split, which we corrected manually.

We also include experiments based on the `rag-mini-bioasq` dataset, henceforth BioASQ, which is a subset of the training dataset that was used for the BioASQ Challenge.⁷ This dataset contains biomedical questions, as well as answers and documents. For the evaluation, we merge its training and test splits, which adds up to 4,719 examples.

5.2 Metrics

We primarily compare each H&S pipeline’s answer with the expected answer from the dataset using several metrics. For some experiments on the RepliQA dataset, we also compare the outputs of the highlighter against the gold passage (`long_answer`). We evaluate the quality of the responses based on three types of evaluation metrics, as summarized in Table 1.

Token-based metrics. These provide hard numbers based on a comparison between the tokens of the expected answer (or gold passage) and the provided answer. We employ two such metrics, based on the work by Adlakha et al. [3], which evaluate the correctness (“Recall”) and faithfulness (“K-Precision”) of an answer.

⁴<https://learn.microsoft.com/en-us/azure/ai-services/openai/how-to/structured-outputs>

⁵<https://github.com/rapidfuzz/RapidFuzz> (threshold = 95)

⁶<https://huggingface.co/deepset/deberta-v3-base-squad2>

⁷<https://www.bioasq.org>

Table 1: Metrics used for evaluating individual H&S components and the full H&S pipeline.

Name	Type	Description	Ref
Recall	Token	Proportion of tokens in the reference answer are present in the model’s response.	[3]
K-Precision	Token	Proportion of tokens in the model’s response that are present in the gold passage.	[3]
Poll Multihop Correctness	LLM	Correctness of a generated response against a reference answer using few-shot learning.	[34]
Reliable CI Relevance	LLM	Relevance of a passage to a query based on a four-point scale: Irrelevant, Related, Highly relevant, Perfectly relevant.	[24]
MTBench Chat Bot Response Quality	LLM	Quality of the response based on helpfulness, relevance, accuracy, depth, creativity, and level of detail, assigning a numerical grade.	[37]
ComparisonJudge	LLM	Compare two answers to the same question and select either a winner, a tie, or a tie where neither answer is acceptable.	Ours

LLM-as-a-Judge. We use prompt-tuned LLMs that rate the responses w.r.t. various criteria [7, 37]. Concretely, we adopt three standard LLM-as-a-Judge implementations, so as to provide a diverse judgment [24, 34, 37].

Comparison judge. We use a zero-shot prompt-tuned LLM with the task of deciding a “winner” between two alternative answers to a question. This judge also has the option to declare a “tie” between the answers, or to decide that neither answer is acceptable. In our implementation, we randomize the order of the two answers to mitigate any potential bias.

6 Evaluating H&S components

In this section we empirically study the individual components of an H&S pipeline. First, we explore whether H&S is needed at all, or whether a simple passage retrieval or extraction Q&A pipeline (i.e., highlighter without the summarizer) may suffice. Next, we compare our different highlighter implementations in terms of their K-Precision and Recall. Finally, we investigate whether the summarizer can recreate the original question based solely on the outputs of a highlighter.

6.1 Do we need a generative summarizer?

First, one may wonder: do we actually need the summarizer in the pipeline? Given the long history of passage retrieval and extractive Q&A, a simpler highlighter, which can be implemented on the basis of an extractive Q&A model, may perform well in Q&A tasks.

To answer this question, we compare the outputs of the individual highlighters against the response given by the full H&S pipeline using a ComparisonJudge. As shown in Figure 2, our results indicate that using a generative LLM (summarizer) after the highlighting step can lead to noticeable improvements. This is particularly evident on the BioASQ dataset, where the H&S outputs were preferred in 58% of cases versus only 6% for the highlighter alone. Based on a qualitative manual inspection of the ComparisonJudge’s explanations accompanying its ratings, we observe that the H&S pipeline is typically preferred for its more “natural” and “well-structured” answers, whereas the highlighter is preferred for more “concise” and “direct quote” answers.

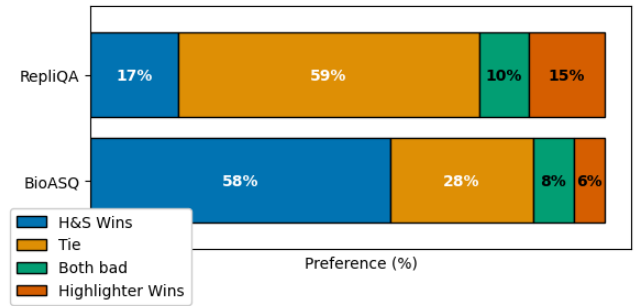


Figure 2: Do we need H&S, or can we just do highlighting (e.g., extractive Q&A)? The preference is based on an LLM-based ComparisonJudge.

6.2 How good are the highlighters?

We now compare the various highlighter implementations, in terms of their K-Precision and Recall with respect to the gold passage. Since this analysis requires a gold passage, we are only able to use the RepliQA dataset.

As shown in Table 2, we observe that LLM-based highlighters significantly outperform those based on DeBERTaV3, especially in terms of Recall. However, as discussed in section 7, LLM-based highlighters incur higher computational overheads, which may be a consideration in practical applications.

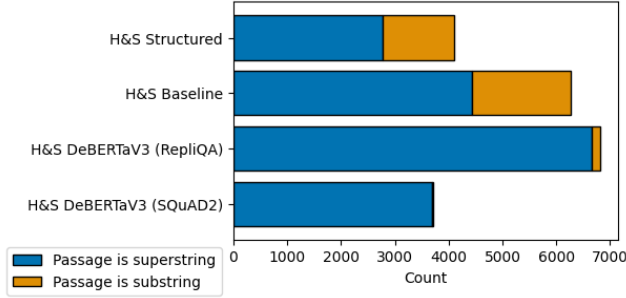
We did not focus on optimizing the performance of the DeBERTaV3 models, and we hypothesize that there may still be scope for further improvement. An encouraging result in this direction is that fine-tuning DeBERTaV3 for the specific task of long context highlighting (RepliQA dataset) improves performance compared to the same model fine-tuned on SQuAD2 (short answers).

Results in Table 2 also indicate that using structured outputs for LLM-based highlighters helps performance (Recall); this, as reported in section 7, comes with additional computational costs.

We further evaluate highlighters on a stricter metric: we count how many times the highlighter’s output is either a substring or a superstring of the gold passage; one should bear in mind that the choice of the human-chosen gold passage is somewhat arbitrary. Figure 3 shows that, in this case, the DeBERTaV3 highlighter that was fine-tuned on RepliQA has the best performance. We attribute

Table 2: Comparison between highlighter implementations with respect to the gold passage (RepliQA dataset.)

Implementation	K-Precision	Recall
H&S Structured	0.84	0.76
H&S Baseline	0.84	0.65
H&S DeBERTaV3 (RepliQA)	0.80	0.36
H&S DeBERTaV3 (SQuAD2)	0.55	0.22

**Figure 3: How many times the highlighter produced a substring or superstring of the human-curated gold passage on the RepliQA dataset.**

this to the fact that, even if its fine-tuning was done on a separate split of the RepliQA dataset, the model may have captured the style of gold passage selection of the human annotators for RepliQA.

6.3 Can the Summarizer guess the question?

To gain a deeper understanding of the internal workings of an H&S pipeline, we investigate whether the summarizer can guess the user’s original question, based solely on the text provided by the highlighter. We start by noting that this problem is, by its nature, ill-posed. For example, consider the passage “*In his later years, Kant lived a strictly ordered life. It was said that neighbors would set their clocks by his daily walks*”⁸. This text could be a plausible answer to several realistic questions, including “*Was Immanuel Kant a creature of habit?*”, “*What philosopher is best known for their punctuality?*”, or “*Did Kant enjoy walking?*”.

When implementing the summarizer (section 4), we ask the LLM to output `guessed_question`, a field that is not used by the H&S pipeline, but which helps with grounding the model. In Table 3 we measure the K-Precision and Recall of `guessed_question` with respect to the true question. These experiments are carried out for H&S Structured pipeline. As expected, the summarizer is rarely able to guess the question correctly. However, its average performance is far from null, with better pipelines achieving higher scores.

In Table 4, we report some of the *worst* examples, which highlight the ill-posedness of the problem. In section 9 we consider whether monitoring the guessed question can help with controlling the performance of H&S pipelines.

⁸https://en.wikipedia.org/wiki/Immanuel_Kant.

Table 3: Can a summarizer guess the question that was asked?

Pipeline	K-Precision	Recall
RepliQA		
H&S Structured	0.55	0.49
H&S Baseline	0.46	0.39
H&S DeBERTaV3 (RepliQA)	0.41	0.34
H&S DeBERTaV3 (SQuAD2)	0.29	0.21
BioASQ		
H&S Structured	0.38	0.60
H&S Baseline	0.34	0.53
H&S DeBERTaV3 (RepliQA)	0.34	0.46
H&S DeBERTaV3 (SQuAD2)	0.24	0.33

Table 4: Worst question guessed by the summarizer of the H&S Structured pipeline on the RepliQA dataset. K-Precision for all was 0.

Real question	Guessed question
What is the main argument of the opinion piece titled ‘The Critical Shield’?	Why should mandatory automatic software updates be implemented on all consumer devices?
Which local candidate became notorious for tweeting vague statements like ‘New policies on the horizon #blessed #nofilter’?	Who is Harvey Peterson and why did he gain notoriety?
What was the primary topic of discussion by Mark Downey at the Online Marketing Summit in San Diego on October 6, 2023?	How can social media benefit small and medium-sized enterprises (SMEs)?

7 Evaluating the full H&S pipeline

Having evaluated the individual components of H&S, we now compare full H&S pipelines with vanilla RAG. Since H&S only modifies generation step in RAG, we hold the retrieval step constant for each comparison; that is, we compare the pipelines on their ability to answer questions on the basis of a given (retrieved) document.

In this section, we compare the pipelines using all three LLM-as-Judges and the ComparisonJudge, as well as on their ability to decline to answer when they believe no answer can be provided. We also compare time measurements for each of the pipelines (Table 7). Although the absolute time measurements are specific to our experimental setup, this shows the relative time overhead of each pipeline.

7.1 Pairwise comparisons

We use the ComparisonJudge to perform pairwise comparisons between all pipelines. A summary of the results is provided in Table 5, which includes the win rate of each pipeline, as well as its Elo score. The Elo score, invented in the context of chess for

Table 5: Average wins and Elo scores for the pipelines evaluated pairwise using a ComparisonJudge. Ties are omitted.

	Wins	Elo
RepliQA		
H&S Structured	59%	1211
H&S Baseline	54%	1147
H&S DeBERTaV3 (RepliQA)	32%	992
Vanilla RAG	32%	958
H&S DeBERTaV3 (SQuAD2)	8%	688
BioASQ		
H&S Structured	69%	1257
H&S Baseline	62%	1109
H&S DeBERTaV3 (RepliQA)	24%	945
Vanilla RAG	42%	899
H&S DeBERTaV3 (SQuAD2)	10%	787

ranking players based on 1-to-1 matches [11], has recently become a popular metric for comparing LLMs (e.g., by Chiang et al. [8]). The full pairwise results are shown in Figure 4 and Table 9.

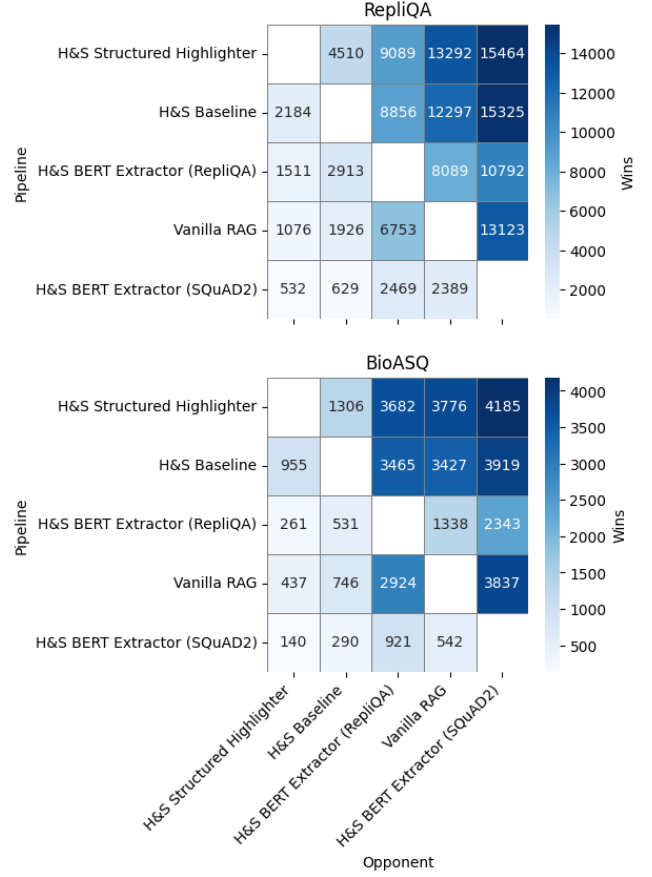
We observe that the LLM-based H&S pipelines (Structured and Baseline) vastly outperform the other pipelines. **This is a surprising result: it suggests that H&S simultaneously gives security against malicious users as well as better response accuracy than vanilla RAG.**

7.2 Graders

We compare the pipelines on the basis of three LLM-as-Judges that independently evaluate correctness, relevance, and quality (as detailed in Table 1).

First, as presented in Figure 5, we count the number of times each pipeline received the highest rating, including ties, out of the tested pipelines (a “win”). We observe that, whereas pipelines are nearly equivalent when it comes to relevance and quality of the responses, the correctness shows some separation. On RepliQA, H&S Structured obtains a better result (13,698 wins out of 17,955 examples), followed by H&S Baseline (13,401), and Vanilla RAG (13,332). We observe a similar behavior on BioASQ. Interestingly, on the BioASQ dataset we note that H&S DeBERTaV3 (SQuAD2) outperforms H&S DeBERTaV3 (RepliQA); this is likely due to the fact that the latter was fine-tuned on (a separate split of) RepliQA. Figure 6 shows the distribution of ratings from each judge. We discuss each judge separately.

Correctness. Poll Multihop Correctness evaluates the correctness of a generated answer against the reference answer. Most pipelines output >80% correct answers. We observe that H&S DeBERTaV3 (SQuAD2) performs poorly on the RepliQA dataset, likely due to its base highlighter being fine-tuned for returning short answers. On the other hand, we observe that the same base model fine-tuned on a separate split of the RepliQA dataset performs rather well on the RepliQA dataset, suggesting that the performance of such extractive Q&A models can be further improved through fine-tuning. However, the same pipeline performs poorly on the BioASQ dataset. This suggests that the data distribution highly impacts

**Figure 4: Wins of each pipeline in a pairwise comparison based on the ComparisonJudge. Ties are omitted.**

performance, and that more varied fine-tuning datasets should be used in practice.

Relevance. The Reliable CI Relevance explicitly captures how relevant an answer is to a question. We observe minimal discrepancy as all pipelines receive a nearly perfect relevance on average (2.72 out of 3). We observe slightly better performance by the vanilla RAG pipeline, followed by the H&S DeBERTaV3 (SQuAD2).

Response quality. MTBench Chat Bot Response Quality evaluates the responses on a scale from 1 to 10. We observe its rating for all pipelines is just below (resp. above) the 5 out of 10 rating for the RepliQA (resp. BioASQ) datasets. Based on manual inspection of the judge’s explanations, we observe that the judge tends to demand additional information beyond what is available in the retrieved document. This suggests the judge is better suited for rating LLMs on open ended questions, without source documents, rather than RAG settings.

7.3 Decline to answer

Around 10% of questions in the RepliQA dataset cannot be answered based on the provided document. We evaluate how well the various pipelines decline to answer these questions.

Highlight, then Summarize: RAG without the jailbreaks

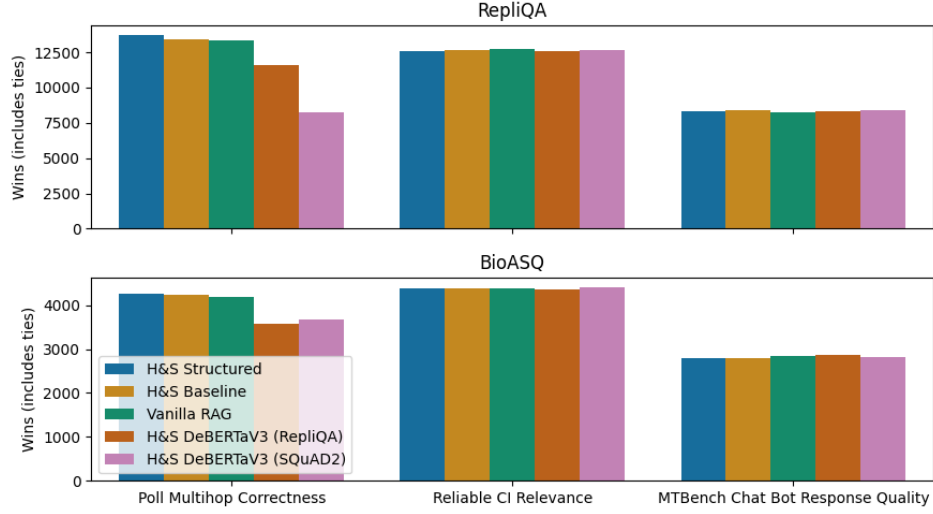


Figure 5: Number of times each pipeline had the highest rating, including ties, out of all pipelines, according to LLM-as-Judges.

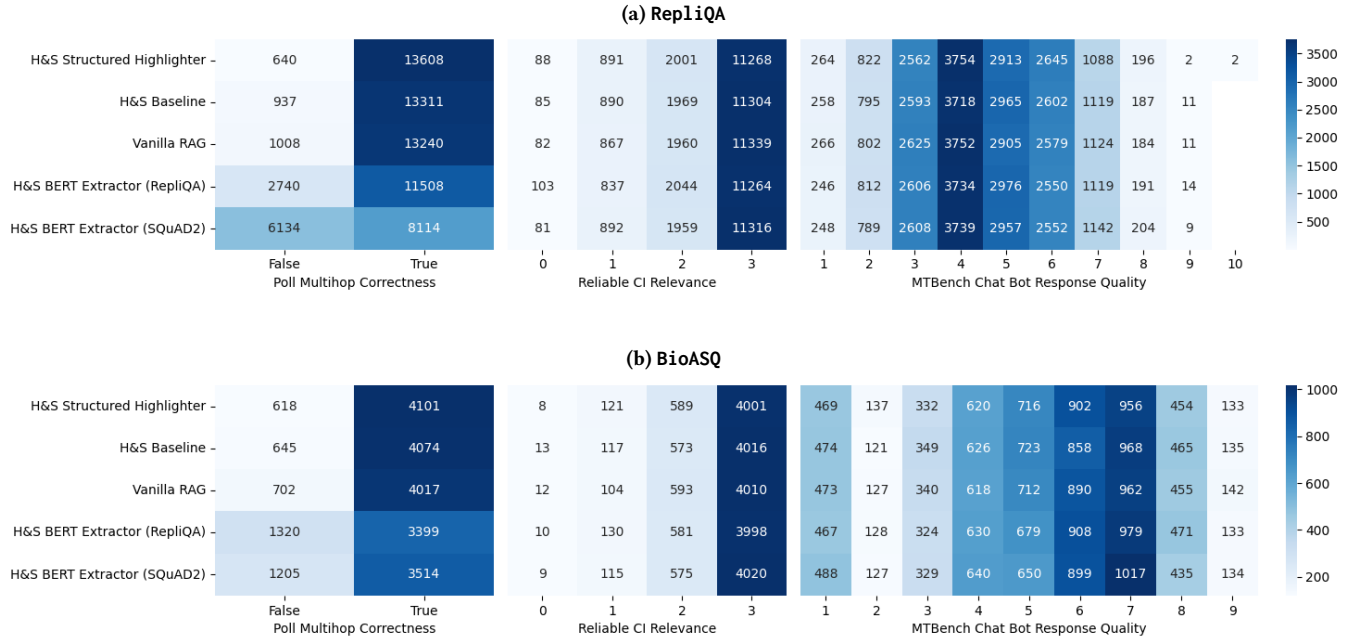


Figure 6: Response’s evaluation via LLM-as-Judges on the two datasets, measuring: correctness, relevance, and quality.

In Table 6, we report K-Precision and Recall for each pipeline in terms of declining to answer. We observe that DeBERTaV3-based models fare well, which is possibly due to the fact that their fine-tuning considers the no-answer case explicitly. The H&S DeBERTaV3 (RepliQA) performs particularly well; this may be partially due to the fact that it was trained on data with a similar distribution as the test set. H&S Structured also does relatively well, possibly because the structured output option helps its reflection process.

Nevertheless, the absolute performance of all the pipelines suggests there is still significant potential for improvement in this aspect. Future H&S implementations should consider either better fine-tuning or few-shot prompt-tuning to improve this metric.

7.4 Processing time

Finally, we measure the time taken by each pipeline to generate an answer, averaged over 40 questions, as shown in Table 7. The

Table 6: Precision and Recall for declined answers.

Pipeline	Precision	Recall	F1
H&S DeBERTaV3 (RepliQA)	0.85	0.99	0.91
H&S Structured	0.94	0.39	0.55
H&S DeBERTaV3 (SQuAD2)	0.55	0.48	0.51
Vanilla RAG	0.95	0.32	0.48
H&S Baseline	0.58	0.20	0.29

absolute time values are specific to our setup, but the relative differences show that, as expected, the LLM-based H&S pipelines require more time to generate an answer. We note that our implementation has not yet been optimized for performance.

Table 7: Average time to process one question, averaged across 40 examples.

Pipeline	Time (s)
QAEvaluator	0.49
H&S DeBERTaV3 (RepliQA)	0.75
H&S DeBERTaV3 (SQuAD2)	0.76
H&S Baseline	2.26
H&S Structured	3.05

8 Security considerations

Fundamentally, all forms of jailbreaking involve some type of adversarial input to an LLM, which potentially causes the LLM to generate undesirable outputs. In the widely-used setting of a RAG-powered system with a trusted knowledge base, H&S mitigates jailbreaks *by design* by ensuring the generative LLM that produces the final output (i.e., the summarizer) never receives direct adversarial inputs, even if these are explicitly provided by the user. In this section, we demonstrate this security property using both a non-adaptive attack and new adaptive attack specifically designed to subvert H&S.

8.1 Non-Adaptive Attack

We begin with a non-adaptive (i.e., generic) attack that aims to jailbreak any LLM-based system.

Data and setup. Evaluating LLM jailbreaks is notoriously challenging because it requires an automated method for detecting whether the jailbreak succeeded. Inspired by the recent LLMail-Inject challenge [1], we use *tool calling* as a well-defined attack target. We give the LLM the ability to call a (simulated) email sending tool by outputting a specific string (e.g., `send_email()`). For example, a customer service chatbot may have this type of tool to send an email to the customer support team if a user asks about certain topics, but the users should not be able to control when the LLM triggers this tool call. The attacker’s goal is therefore to cause the LLM to trigger this tool call. When evaluating H&S, we give *both* LLMs (i.e., highlighter and summarizer) the ability to trigger this tool call. However, in practice only the summarizer would have

Table 8: Security evasion of RAG and H&S pipelines using the LLMail-Inject competition dataset (Scenario 2). Evasion of the Highlighter *does not* constitute a bypass of the H&S pipeline, and it is only reported for comparison.

	Tool called	Arguments valid
RAG	81%	53%
H&S (Highlighter only)	93%	63%
H&S (Full)	0%	0%

this ability, so a successful attack against H&S would need to trick the summarizer into making this tool call.

For these experiments, we use 1,028 successful attack prompts from the LLMail-Inject challenge dataset (Scenario 2) [1]. Our setup differs slightly from that of the challenge: in the challenge, the adversarial inputs were retrieved from the knowledge base, whereas in our setting we input these directly as the user’s prompt. However, this modified setup does not affect the results: successful prompts from LLMail-Inject are usually successful in our setup.

Results. Table 8 shows the percentage of inputs for which the tool was called (with any arguments) or called with valid arguments. As expected, the vast majority of the attack inputs succeed in jailbreaking the RAG pipeline, but none are successful against H&S. The results on H&S (Highlighter only) confirm that our highlighter LLM, like any other, is still susceptible to jailbreaks when processing direct adversarial input. This further underscores the need for system-level defenses, such as H&S.

8.2 Adaptive H&S-specific Attack

As explained in section 3, we assume a strong adversary who may have full knowledge of H&S and the contents of the knowledge base used by the system under attack. We envision that such an attacker could attempt to subvert H&S with the following adaptive attack, which is somewhat reminiscent of traditional system attacks such as Return-Oriented Programming (ROP) [33]. For a target sentence, the attacker i) browses the RAG knowledge base to find a document that contains (most of) the words in the target sentence, ii) causes that document to be retrieved, and iii) asks the H&S highlighter to only highlight the desired words.

For example, suppose the attacker wishes the system to output “You won a \$10 voucher”. If the attacker can persuade the H&S highlighter to highlight only the bold words from the retrieved document shown in Figure 7, they will likely succeed in having the summarizer output their desired response.

While this attack is theoretically possible, its difficulty can be substantially increased through a simple H&S design constraint: we only allow the highlighter to highlight contiguous passages of at least a certain length. In our implementations, this attack was readily prevented by the use of fuzzy string matching to ensure that the highlighted text corresponds to a contiguous portion of the original text. Fundamentally, H&S transforms the very challenging problem of detecting any possible jailbreak in the LLM’s inputs into the significantly simpler task of inspecting the knowledge base to check that it does not contain strings that can be used as gadgets.

Theoretical gadgets-based attack against H&S

You may submit a request for reimbursement if the expense falls within the approved categories outlined in Appendix D. Once the claim is verified and approved, it is recorded as a completed transaction and marked as “won” in the internal tracking system. A confirmation email will be sent to the claimant within three (3) business days. If the reimbursement includes a \$10 threshold adjustment, the system will automatically generate a voucher code for accounting purposes. This **voucher** is non-transferable and must be used within the same fiscal quarter.

Figure 7: Theoretical gadget-based attack against H&S

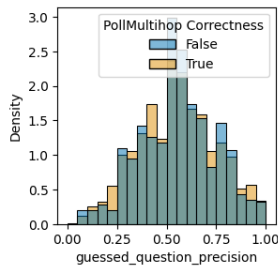


Figure 8: Lack of correlation between the ability of the summarizer to guess the question and the performance of the H&S Structured pipeline. Pearson correlation: -0.004.

9 Discussion and future directions

This study of a new design pattern opens up a large number of research questions and directions to explore. We discuss some of the challenges that we leave open for future work.

Using guessed question as grounding data. Future work or implementations of H&S pipelines could monitor the distance between the actual question and the guessed questions: if too far, the system could be programmed to refuse to answer, or to adapt the summarizer’s task accordingly. From a security perspective, we note that this would have no impact: as long as the summarizer does not see the actual input from the user, the system can still be considered to be safe from jailbreaking.

In our current implementation, we found no significant correlation between the pipeline’s performance and the summarizer’s ability to predict the question. For example, Figure 8 shows the relationship between the full pipeline’s performance (measured via an LLM as a judge) and the K-Precision of the summarizer at guessing the question.

Automated inspection for gadgets. Our threat model assumes a trustworthy knowledge base. While this is realistic in many real-world use cases (e.g., Q&A based on FAQs), it may not be the case

in all RAG applications. We observe, however, that H&S makes it easier to look for threats: it reduces the scope of the problem from inspecting *both* the knowledge base and the very large space of possible user inputs, down to inspecting only the knowledge base. We also note that this makes it much easier to evaluate the security of a deployment. For example, one can scan their knowledge base for potential triggers, by inspecting each document via a (sequential) sliding window representative of what the highlighter is allowed to highlight. We encourage future research to develop methods in this direction.

Handling Yes/No questions. In some cases, a simple “Yes” or “No” can be a satisfactory answer. We observe that an H&S system can be easily augmented to provide a better user experience in these circumstances, whilst ensuring the same level of security. One example is as follows: The highlighter can optionally pass a tag to the summarizer indicating: i) whether the question being asked is of yes-no type, and if so ii) what it thinks to be the answer (yes/no); since both of these fields are of boolean type, we expect them not to increase the system’s risk. The summarizer can then be instructed via prompt-tuning to augment its question with this information.

Handling multiple questions. In some cases, the user of a RAG pipeline may be asking more than one question in the same prompt. We suggest this can be tackled via system design. For example, the highlighter can first split the questions that are present in the prompt, and then have the H&S system make separate calls to the summarizer. We leave the evaluation of this setup to future work.

Highlighting non-contiguous parts of the document. LLM-based and BERT-based highlighters are readily able to highlight potentially non-contiguous parts of a document, if necessary. Future work should study this more in depth, keeping in mind the risks described in subsection 8.2.

Constrained decoding. A prompt-tuned LLM highlighter can be seemingly improved via constrained decoding: via this process, the LLM is forced to match against the original text, thereby reducing hallucinations.⁹ Future work may consider its use for building highlighters in the case of H&S pipelines.

Does H&S reduce hallucinations? Intuitively, passage retrieval and extractive Q&A should mitigate hallucinations, since quoting the original document is arguably a better controllable process than LLM-based generative Q&A. However, it is possible that the highlighter selects portions from the document that are either irrelevant or misleading; this was observed in section 7, where H&S (as well as vanilla RAG) often tried to answer unanswerable questions. Furthermore, the summarizer has uncertainty about the question, which may be kept in mind when building H&S systems; however, in practice, this did not seem to negatively affect the performance. We encourage further empirical studies into hallucinations for H&S.

10 Conclusions

We proposed Highlight&Summarize (H&S), a design pattern that enhances the generative part of a RAG pipeline to provide security by design against malicious users. Our empirical evaluations of the

⁹This was illustrated, for example, in <https://yonigottesman.github.io/2023/08/10/extractive-generative.html>.

performance and security of H&S demonstrate that this approach can both improve the accuracy of responses, compared to vanilla RAG, whilst virtually eliminating jailbreaks from malicious users. We also discussed several new research questions that arise from this design pattern, and encourage their study as exciting avenues for future research.

Acknowledgments

We are grateful to Sahar Abdelnabi and Santiago Zanella-Beguelin for useful discussions. We thank Elisa Alessandrini for help with revising this manuscript.

References

- [1] Sahar Abdelnabi, Aileen Fay, Ahmed Salem, Egor Zverev, Kai-Chieh Liao, Chi-Huang Liu, Chun-Chih Kuo, Jannis Weigend, Danyael Manlangit, Alex Apostolov, et al. 2025. LLMail-Inject: A Dataset from a Realistic Adaptive Prompt Injection Challenge. *arXiv preprint arXiv:2506.09956* (2025).
- [2] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).
- [3] Vaibhav Adlakha, Parishad BehnamGhader, Xing Han Lu, Nicholas Meade, and Siva Reddy. 2024. Evaluating correctness and faithfulness of instruction-following models for question answering. *Transactions of the Association for Computational Linguistics* 12 (2024), 681–699.
- [4] Gabriel Alon and Michael Kamfonas. 2023. Detecting Language Model Attacks with Perplexity. *arXiv:2308.14132* [cs.CL] <https://arxiv.org/abs/2308.14132>
- [5] Luca Beurer-Kellner, Beat Buesser, Ana-Maria Cretu, Edoardo Debenedetti, Daniel Dobos, Daniel Fabian, Marc Fischer, David Froelicher, Kathrin Grosse, Daniel Naef, Ezinwanne Ozoani, Andrew Paverd, Florian Tramèr, and Václav Volhejn. 2025. Design Patterns for Securing LLM Agents against Prompt Injections. *arXiv:2506.08837* [cs.LG] <https://arxiv.org/abs/2506.08837>
- [6] James P. Callan. 1994. Passage-level evidence in document retrieval. In *SIGIR '94: Proceedings of the Seventeenth Annual International ACM-SIGIR Conference on Research and Development in Information Retrieval, organised by Dublin City University*. Springer, 302–310.
- [7] Cheng-Han Chiang and Hung-yi Lee. 2023. Can large language models be an alternative to human evaluations? *arXiv preprint arXiv:2305.01937* (2023).
- [8] Wei-Lin Chiang, Lianmin Zheng, Ying Sheng, Anastasios Nikolas Angelopoulos, Tianle Li, Dacheng Li, Banghua Zhu, Hao Zhang, Michael Jordan, Joseph E. Gonzalez, et al. 2024. Chatbot arena: An open platform for evaluating llms by human preference. In *Forty-first International Conference on Machine Learning*.
- [9] Manuel Costa, Boris Köpf, Aashish Kolluri, Andrew Paverd, Mark Russinovich, Ahmed Salem, Shruti Tople, Lukas Wutschitz, and Santiago Zanella-Béguelin. 2025. Securing AI Agents with Information-Flow Control. *arXiv:2505.23643* [cs.CR] <https://arxiv.org/abs/2505.23643>
- [10] Edoardo Debenedetti, Ilia Shumailov, Tianqi Fan, Jamie Hayes, Nicholas Carlini, Daniel Fabian, Christoph Kern, Chongyang Shi, Andreas Terzis, and Florian Tramèr. 2025. Defeating Prompt Injections by Design. *arXiv:2503.18813* [cs.CR] <https://arxiv.org/abs/2503.18813>
- [11] Arpad E. Elo. 1967. The proposed uscf rating system, its development, theory, and applications. *Chess life* 22, 8 (1967), 242–247.
- [12] Zhengmian Hu, Gang Wu, Saayan Mitra, Ruiyi Zhang, Tong Sun, Heng Huang, and Viswanathan Swaminathan. 2024. Token-Level Adversarial Prompt Detection Based on Perplexity Measures and Contextual Information. *arXiv:2311.11509* [cs.CL] <https://arxiv.org/abs/2311.11509>
- [13] Neel Jain, Avi Schwarzschild, Yuxin Wen, Gowthami Somepalli, John Kirchenbauer, Ping yeh Chiang, Micah Goldblum, Aniruddha Saha, Jonas Geiping, and Tom Goldstein. 2023. Baseline Defenses for Adversarial Attacks Against Aligned Language Models. *arXiv:2309.00614* [cs.LG] <https://arxiv.org/abs/2309.00614>
- [14] Changyue Jiang, Xudong Pan, Geng Hong, Chenfu Bao, and Min Yang. 2024. RAG-Thief: Scalable Extraction of Private Data from Retrieval-Augmented Generation Applications with Agent-based Attacks. *arXiv:2411.14110* [cs.CR] <https://arxiv.org/abs/2411.14110>
- [15] Jing Jiang and Chengxiang Zhai. 2006. Extraction of coherent relevant passages using hidden markov models. *ACM Transactions on Information Systems (TOIS)* 24, 3 (2006), 295–319.
- [16] Daniel Jurafsky and James H. Martin. 2025. *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition with Language Models* (3rd ed.). <https://web.stanford.edu/~jurafsky/slp3/> Online manuscript released January 12, 2025.
- [17] Aounon Kumar, Chirag Agarwal, Suraj Srinivas, Aaron Jiaxun Li, Soheil Feizi, and Himabindu Lakkaraju. 2025. Certifying LLM Safety against Adversarial Prompting. *arXiv:2309.02705* [cs.CL] <https://arxiv.org/abs/2309.02705>
- [18] Kenton Lee, Shimi Salant, Tom Kwiatkowski, Ankur Parikh, Dipanjan Das, and Jonathan Berant. 2016. Learning recurrent span representations for extractive question answering. *arXiv preprint arXiv:1611.01436* (2016).
- [19] Microsoft. 2025. Safety system messages. <https://learn.microsoft.com/en-us/azure/ai-services/openai/concepts/system-message>
- [20] Joao Monteiro, Pierre-Andre Noel, Etienne Marcotte, Sai Rajeswar Mudumba, Valentina Zantedeschi, David Vazquez, Nicolas Chapados, Chris Pal, and Perouz Taslakian. 2024. RepliQA: A question-answering dataset for benchmarking llms on unseen reference content. *Advances in Neural Information Processing Systems* 37 (2024), 24242–24276.
- [21] Norman Mu, Jonathan Lu, Michael Lavery, and David Wagner. 2025. A Closer Look at System Prompt Robustness. *arXiv:2502.12197* [cs.CL] <https://arxiv.org/abs/2502.12197>
- [22] John O'Connor. 1975. Retrieval of answer-sentences and answer-figures from papers by text searching. *Information Processing & Management* 11, 5-7 (1975), 155–164.
- [23] John O'Connor. 1980. Answer-passage retrieval by text searching. *Journal of the American Society for Information Science* 31, 4 (1980), 227–239.
- [24] Harrie Oosterhuis, Rolf Jagerman, Zhen Qin, Xuanhui Wang, and Michael Bender-sky. 2024. Reliable confidence intervals for information retrieval evaluation using generative ai. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*.
- [25] OpenAI. [n.d.]. Moderation. <https://platform.openai.com/docs/guides/moderation/overview>
- [26] Kate Pearce, Tiffany Zhan, Aneesh Komanduri, and Justin Zhan. 2021. A comparative study of transformer-based language models on extractive question answering. *arXiv preprint arXiv:2110.03142* (2021).
- [27] Yuefeng Peng, Junda Wang, Hong Yu, and Amir Houmansadr. 2025. Data Extraction Attacks in Retrieval-Augmented Generation via Backdoors. *arXiv:2411.01705* [cs.CR] <https://arxiv.org/abs/2411.01705>
- [28] Matthew Pisano, Peter Ly, Abraham Sanders, Bingsheng Yao, Dakuo Wang, Tomek Strzalkowski, and Mei Si. 2024. Bergeron: Combating Adversarial Attacks through a Conscience-Based Alignment Framework. *arXiv:2312.00029* [cs.CR] <https://arxiv.org/abs/2312.00029>
- [29] Archiki Prasad, Trung Bui, Seunghyun Yoon, Hanieh Deilamsalehy, Franck Der-noncourt, and Mohit Bansal. 2023. MeetingQA: Extractive question-answering on meeting transcripts. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 15000–15025.
- [30] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ Questions for Machine Comprehension of Text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, Jian Su, Kevin Duh, and Xavier Carreras (Eds.). Association for Computational Linguistics, Austin, Texas, 2383–2392. doi:10.18653/v1/D16-1264 *arXiv:1606.05250* [cs.CL]
- [31] Traian Rebedea, Razvan Dinu, Makesh Sreedhar, Christopher Parisien, and Jonathan Cohen. 2023. NeMo Guardrails: A Toolkit for Controllable and Safe LLM Applications with Programmable Rails. *arXiv:2310.10501* [cs.CL] <https://arxiv.org/abs/2310.10501>
- [32] Alexander Robey, Eric Wong, Hamed Hassani, and George J. Pappas. 2024. SmoothLLM: Defending Large Language Models Against Jailbreaking Attacks. *arXiv:2310.03684* [cs.LG] <https://arxiv.org/abs/2310.03684>
- [33] Hovav Shacham. 2007. The geometry of innocent flesh on the bone: return-into-libc without function calls (on the x86). In *Proceedings of the 14th ACM Conference on Computer and Communications Security (CCS '07)*. doi:10.1145/1315245.1315313
- [34] Pat Verga, Sebastian Hofstatter, Sophia Althammer, Yixuan Su, Aleksandra Piktus, Arkady Arkhangorodsky, Minjie Xu, Naomi White, and Patrick Lewis. 2024. Replacing judges with juries: Evaluating LLM generations with a panel of diverse models. *arXiv preprint arXiv:2404.18796* (2024).
- [35] Luqi Wang, Kaiwen Zheng, Liyin Qian, and Sheng Li. 2022. A survey of extractive question answering. In *2022 International Conference on High Performance Big Data and Intelligent Systems (HDIS)*. IEEE, 147–153.
- [36] Baolei Zhang, Yuxi Chen, Minghong Fang, Zhuqing Liu, Lihai Nie, Tong Li, and Zheli Liu. 2025. Practical Poisoning Attacks against Retrieval-Augmented Generation. *arXiv:2504.03957* [cs.CR] <https://arxiv.org/abs/2504.03957>
- [37] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric Xing, et al. 2023. Judging llm-as-a-judge with mt-bench and chatbot arena. *Advances in Neural Information Processing Systems* (2023).
- [38] Wei Zou, Rumpeng Geng, Binghui Wang, and Jinyuan Jia. 2024. PoisonedRAG: Knowledge Corruption Attacks to Retrieval-Augmented Generation of Large Language Models. *arXiv:2402.07867* [cs.CR] <https://arxiv.org/abs/2402.07867>

A Additional results

Table 9: Pairwise battle results: Each cell shows (Wins, Ties) for the row model vs the column model. The last columns show total wins and ties.

RepliQA Opponent	RAG		DeBERTaV3 (SQuAD2)		Baseline		Structured		DeBERTaV3 (RepliQA)		Total	
	Wins	Ties	Wins	Ties	Wins	Ties	Wins	Ties	Wins	Ties	Wins	Ties
RAG	–	–	13123	2443	1076	3587	1926	3732	6753	3113	22878	12875
DeBERTaV3 (SQuAD2)	2389	2443	–	–	532	1959	629	2001	2469	4694	6019	11097
Structured	13292	3587	15464	1959	–	–	4510	11261	9089	7355	42355	24162
Baseline	12297	3732	15325	2001	2184	11261	–	–	8856	6186	38662	23180
DeBERTaV3 (RepliQA)	8089	3113	10792	4694	1511	7355	2913	6186	–	–	23305	21348

BioASQ Opponent	RAG		DeBERTaV3 (SQuAD2)		Baseline		Structured		DeBERTaV3 (RepliQA)		Total	
	Wins	Ties	Wins	Ties	Wins	Ties	Wins	Ties	Wins	Ties	Wins	Ties
RAG	–	–	3837	340	746	546	437	506	2924	457	7944	1849
DeBERTaV3 (SQuAD2)	542	340	–	–	290	510	140	394	921	1455	1893	2699
Baseline	3427	546	3919	510	–	–	955	2458	3465	723	11766	4237
Structured	3776	506	4185	394	1306	2458	–	–	3682	776	12949	4134
DeBERTaV3 (RepliQA)	1338	457	2343	1455	531	723	261	776	–	–	4473	3411